

Improving accuracy and speeding up Document Image Classification through parallel systems

Javier Ferrando¹, Juan Luis Domínguez¹, Jordi Torres^{1,2}, Raúl García¹,
David García¹, Daniel Garrido³, Jordi Cortada⁴, and Mateo Valero^{1,2}

¹ Barcelona Supercomputing Center - Centro Nacional de Supercomputación
{javier.ferrando,juan.dominguez,jordi.torres,raul.garcia,
david.garcia2,mateo.valero}@bsc.es

² Universitat Politècnica de Catalunya, UPC-BarcelonaTech

³ Serimag Media - TAAD
daniel.garrido@serimagmedia.com

⁴ CaixaBank
jorge.cortada@caixabank.com

Abstract. This paper presents a study showing the benefits of the EfficientNet models compared with heavier Convolutional Neural Networks (CNNs) in the Document Classification task, essential problem in the digitalization process of institutions. We show in the RVL-CDIP dataset that we can improve previous results with a much lighter model and present its transfer learning capabilities on a smaller in-domain dataset such as Tobacco3482. Moreover, we present an ensemble pipeline which is able to boost solely image input by combining image model predictions with the ones generated by BERT model on extracted text by OCR. We also show that the batch size can be effectively increased without hindering its accuracy so that the training process can be sped up by parallelizing throughout multiple GPUs, decreasing the computational time needed. Lastly, we expose the training performance differences between PyTorch and Tensorflow Deep Learning frameworks.

Keywords: Document Image Classification · Deep Learning · Parallel Systems · EfficientNet · BERT · Scalability · TensorFlow · PyTorch.

1 Introduction

Document digitization has become a common practice in a wide variety of industries that deal with vast amounts of archives. Document classification is a task to face when trying to automate their document processes, but high intra-class and low inter-class variability between documents have made this a challenging problem.

First attempts focused on structural similarity between documents [40] and on feature extraction [24,12,30] to differentiate characteristics of each class. The combination of both approaches has also been tested [14].

Several classic machine learning techniques have been applied to these problem, i. e. K-Nearest Neighbor approach [7], Hidden Markov Model [19] and Random Forest Classifier [29,24] while using SURF local descriptors before the Convolutional Neural Networks (CNNs) came into scene.

With the rise of Deep Learning, researchers have tried deep neural networks to improve the accuracy of their classifiers. CNNs have been proposed in past works, initially in 2014 by Le Kang *et al.* [26] who started with a simple 4-layer CNN trained from scratch. Then, transfer learning was demonstrated to work effectively [21,1] by using a network pre-trained on ImageNet [17]. And latest models have become increasingly heavier (greater number of parameters) [46,16,2] as shown in Table 1, with the speed and computational resources drawback this entails.

Recently, textual information has been used by itself or as a combination together with visual features extracted by the previously mentioned models. Although Optical Character Recognition (OCR) is prone to errors, particularly when dealing with handwritten documents, the use of modern Natural Language Processing (NLP) techniques have demonstrated a boost in the classifiers performance [35,6,5].

The contributions of this paper can be summarized in two main topics:

- Algorithmic performance: we propose a model and a training procedure to deal with images and text that outperforms the state-of-the-art in several settings and is lighter than any previous neural network used to classify the BigTobacco dataset, the most popular benchmark for Document Image Classification (Table 1).
- Training process speed up: we demonstrate the ability of these models to maintain their performance while saving a large amount of time by parallelizing over several GPUs. We also show the performance differences between the two most popular Deep Learning frameworks (TensorFlow and Pytorch), when using their own libraries dedicated to this task.

2 Document Image Classification

Document Image Classification task tries to predict the class which a document belongs to by means of analyzing its image representation. This challenge can be tackled in two ways, as an image classification problem and as a text classification problem. The former tries to look for patterns in the pixels of the image to find elements such as shapes or textures that can be associated to a certain class. The latter tries to understand the language written in the document and relate this to the different classes.

2.1 Datasets

As mentioned earlier, in this work we make use of two publicly available datasets containing samples of images from scanned documents from USA Tobacco companies, published by Legacy Tobacco Industry Documents and created by the University of California San Francisco (UCSF). We find these datasets a good representation of what enterprises and institutions may face with, based on the quality and type of classes. Furthermore, they have been go-to datasets in this research field since 2014 with which we can compare results.

RVL-CDIP (Ryerson Vision Lab Complex Document Information Processing) is a 400.000 document sample (BigTobacco from now onwards) presented in [21] for document classification tasks. This dataset contains the first page of each of the documents, which are labeled in 16 different classes with equal number of elements per class. A smaller sample containing 3482 images was proposed in [24] as Tobacco3482 (SmallTobacco henceforth). This dataset is formed by documents belonging to 10 classes not uniformly distributed.

Table 1: Parameters of the CNNs architectures used in BigTobacco.

| Model | #Params |
|-----------------|-------------|
| AlexNet | 60.97M |
| VGG-16 | 138.36M |
| ResNet-50 | 25.56M |
| Inception-V3 | 23.83M |
| EfficientNet-B2 | 9.2M |
| EfficientNet-B0 | 5.3M |

2.2 Deep Learning

The proposed methods in this work are based on supervised Deep Learning, where each document is associated to a class (label) so that the algorithms are trained by minimizing the error between the predictions and the truth. Deep Learning is a branch of machine learning that deals with deep neural networks, where each of the layers is trained to extract higher level representations of the previous ones. These models are trained by solving iteratively an unconstrained optimization problem. In each iteration, a random batch of the training data is fed into the model to compute the loss function value. Then, the gradient of the loss function with respect to the weights of the network is computed (backpropagation) and an update of the weights in the negative direction of the gradient is done. These networks are trained until they converge into a loss function minimum.

2.3 Computer Vision

The field where machines try to get an understanding of visual data is known as Computer Vision (CV). One of the most well-known tasks in CV is image classification. In 2010 The ImageNet Large Scale Visual Recognition Challenge (ILSVRC) was introduced, a competition that dealt with a 1.2 million images dataset belonging to 1000 classes. In 2012 the first CNN-based model significantly reduced the error rate, setting the beginning of the explosion of deep neural networks. From then onwards, deeper networks have become the norm.

The most used architecture in Computer Vision have been CNN-based networks. Their main operation is the convolution one, which consists on a succession of dot products between the vector representations of both the input space ($L_q \times B_q \times d_q$) and the filters ($F_q \times F_q \times d_q$). We slide each filter around the input volume getting an *activation map* of dimension $L_{q+1} = (L_q - F_q + 1)$ and $B_{q+1} =$

$(B_q - F_q + 1)$. The output volume then has a dimension of $L_{q+1} \times B_{q+1} \times d_{q+1}$, where d_{q+1} refers to the number of filters used. We refer to [3] (we used the same notation for simplicity) to a more detailed explanation. Usually, each convolution layer is associated to an activation layer, where an activation function is applied to the whole output volume. To reduce the number of parameters of the network, a pooling layer is typically located between convolution operations. The pooling layer takes a region $P_q \times P_q$ in each of the d_q activation maps and performs an arithmetic operation. The most used pooling layer is the max-pool, which returns the maximum value of the aforementioned region.

2.4 Natural Language Processing

The features learned from the OCR output are achieved by means of Natural Language Processing techniques. NLP is the field that deals with the understanding of human language by computers, which captures underlying meanings and relationships between words.

The way machines deal with words is by means of a real values vector representation. Word2Vec [34] showed that a vector could represent semantic and syntactic relationships between words. CoVe [32] introduced the concept of context-based embeddings, where the same word can have a different vector representation depending on the surrounding text. ELMo [36] followed Cove but with a different training approach, by predicting the next word in a text sequence (Language Modelling), which made it possible to train on large available text corpus. Depending on the task (such as text classification, named entity recognition...) the output of the model can be treated in different ways. Moreover, custom layers can be added to the features extracted by these NLP models. For instance, ULM-Fit [23] introduced a language model and a fine-tuning strategy to effectively adapt the model to various downstream tasks, which pushed transfer learning in the NLP field. Lately, the Transformer architecture [47] has dominated the scene, being the bidirectional Transformer encoder (BERT) [18] the one who established recently state-of-the-art results over several downstream tasks.

3 Related Work

Several ways of measuring models have been shown in the past years regarding document classification on the Legacy Tobacco Industry Documents [31]. Some authors have tested their models on a large-scale sample BigTobacco. Others tried on a smaller version named SmallTobacco, which could be seen as a more realistic scale of annotated data that users might be able to find. Lastly, transfer learning from in-domain datasets has been tested by using BigTobacco to pre-train the models to finally fine-tune on SmallTobacco. Table 2 summarizes the results of previous works in the different categories over time.

First results in the Deep Learning era have been mainly based on CNNs using transfer learning techniques. Multiple networks were trained on specific

sections of the documents [21] to learn region-based high dimensional features later compressed via Principal Component Analysis (PCA). The use of multiple Deep Learning models was also exploited by Arindam Das *et al.* by using an ensemble as a meta-classifier [16]. A VGG-16[41] stack of networks using 5 different classifiers has been proposed, one of them trained on the full document and the others specifically over the header, footer, left body and right body. The Multi Layer Perceptron (MLP) was the ensemble that performed the better. A committee of models but with a SVM as the ensemble was also proposed [37].

Table 2: Previous results comparison (accuracy in %).

| Author | BigTobacco | SmallTobacco | | | |
|-----------------------------|--------------|-------------------------|-------------------|-----------------|--------------|
| | Image | BigTobacco Pre-training | | No Pre-training | |
| | Image | Image | Image + Text | Image | Image + Text |
| Kumar et al. (2014)[24] | | | | 43.8 | |
| Kang et al. (2014)[26] | | | | 65.37 | |
| Afzal et al. (2015)[1] | | | | 77.6 | |
| Harley et al. (2015)[21] | 89.8 | | | 79.9 | |
| Csurka et al. (2016)[15] | 90.7 | | | | |
| Noce et al. (2016)[35] | | | | | 79.8 |
| Afzal et al. (2017)[2] | 90.97 | 91.13 | | | |
| Tensmeyer et al. (2018)[46] | 90.8 | | | | |
| Das et al. (2018)[16] | 92.21 | | | | |
| Audebert et al. (2019)[6] | | | | 84.5 | 87.8 |
| Asim et al. (2019)[5] | | 93.2 ⁵ | 95.8 ⁶ | | |
| Proposed work (2020) | 92.31 | 94.04 | 94.9 | 85.99 | 89.47 |

The addition of content-based information has been investigated on SmallTobacco by extracting text through OCR and embedding the obtained features into the original document images as a previous phase to the training process [35]. Lately, a MobilenetV2 architecture [38] together with a CNN 2D [27,49] taking as input FastText embeddings [9,25] have achieved the best results in SmallTobacco [6].

A study of several CNNs was carried out [2], where VGG-16 architecture was found optimal. Afzal *et al.* also demonstrated that transfer learning from in-domain dataset like BigTobacco increases by a large margin the results in SmallTobacco. This was further investigated by adding content-based information with CNN 2D with ranking textual features (ACC2) to the OCR extracted.

As far as we are concerned, there is no study about the use of multiple GPUs in the training process for the task of Document Image Classification. However, parallelizing a computer vision task has been shown to work properly using ResNet-50, which is a widely used network that usually gives good results despite its low complexity architecture. Several training procedures are demonstrated to

⁵ Accuracy obtained in 9 classes that overlap in BigTobacco

⁶ Evaluation method not specified

work effectively with this model [4,20]. A learning rate value proportional to the batch size, warmup learning rate behaviour, batch normalization, SGD to RMSProp optimizer transition are some of the techniques exposed in these works. A study of the distributed training methods using ResNet-50 architecture on a HPC cluster is shown in [10,11]. To know more about the algorithms used in this field we refer to [8].

4 Proposed Approach

In this section we present the models used and a brief explanation of them. We also show the training procedure used in both BigTobacco and SmallTobacco and the pipeline of our approach to the problem.

4.1 Image model

EfficientNets [45] are a set of light CNNs designed to scale up in a structured manner. The network’s width (w), depth (d) and resolution (r) are defined as: $w = \alpha^\phi$, $d = \beta^\phi$ and $r = \gamma^\phi$, where ϕ is the scaling compound coefficient. The optimization problem is set by constraining $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$ and $\alpha \geq 1, \beta \geq 1, \gamma \geq 1$.

By means of a grid search of α, β, γ with AutoML MNAS framework [44] and fixing $\phi = 1$, a baseline model (B0) is generated optimizing FLOPs and accuracy. Then, the baseline network is scaled up uniformly fixing α, β, γ and increasing ϕ . We find that scaling the resolution parameter as proposed in [45] does not improve the accuracy obtained. In our experiments in Section 5 we proceed with an input image size of 384×384 , which corresponds to a resolution $r = 1.71$, as proposed by Tensmeyer *et al.* in [46] with AlexNet architecture [28].

The main block of the EfficientNets is the mobile inverted bottleneck convolution [38,44]. This block is formed by two linear bottlenecks connected through both a shortcut connection and an intermediate expansion layer with a depth-wise separable convolution (3×3) [13]. Probabilities $P(class|FC)$ are obtained by applying the softmax function on top of the fully connected layer FC of the EfficientNet model.

Pre-training on BigTobacco We train EfficientNets (pre-trained previously on ImageNet) on BigTobacco using Stochastic Gradient Descent for 20 epochs with *Learning Rate Warmup* strategy [22], specifically we follow *STLR* (Slanted Triangular Learning Rate) [23] which linearly increases the learning rate at the beginning of the training process and linearly decreases it after a certain number of iterations. We chose the *reference learning rate* η following the formula proposed in [20] and used in [4] and [22]. Specifically, we set $\eta = 0.2 \cdot \frac{nk}{256}$, where k denotes the number of workers (GPUs) and n the number of samples per worker. Figure 1 shows the multi-GPU training procedure to get EfficientNet_{BigTobacco}, which represents EfficientNet model pre-trained on BigTobacco. EfficientNet is loaded with ImageNet weights (EfficientNet_{ImageNet}) and then located in different GPUs within the same node.

Fine-tuning on SmallTobacco We fine-tune on SmallTobacco the pre-trained models by freezing the entire network but the last softmax layer. Just 5 epochs are enough to get the peak of accuracy. *STLR* is used this time with $\eta = 0.8 \cdot \frac{nk}{256}$. Since only the last layer is trained, we reduce the risk of catastrophic forgetting [33]. Final fine-tuned model is represented as $\text{EfficientNet}_{\text{BigTobacco}}$ in Figure 1.

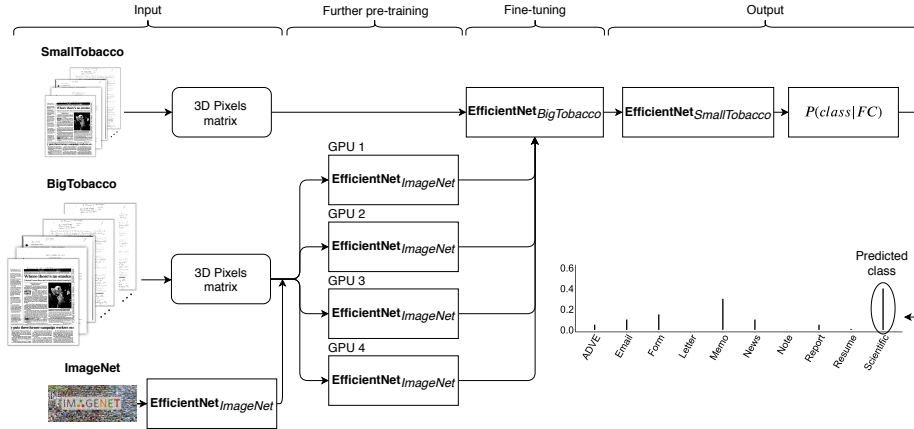


Fig. 1: Pipeline of the different stages of the pre-training of EfficientNet over multiple GPUs.

4.2 Text model

Predictions from OCR Tesseract [42] are obtained by means of the BERT model [18]. BERT is a multi-layer bidirectional Transformer encoder model pre-trained on a large corpus. In this work we use a modification of the original pre-trained $\text{BERT}_{\text{BASE}}$ version. In our case, we reduce to 6 the number of BERT layers since we find less variance in the final results and faster training/inference times. The output vector size is kept to 768. The maximum length of the input sequence is set to 512 tokens. The first token of the sequence is defined as $[CLS]$, while $[SEP]$ is the token used at the end of each sequence.

A fully connected layer is added to the final hidden state of the $[CLS]$ token $h_{[CLS]}$ of the BERT model, which is a representation of the whole sequence. Then, a softmax operation is performed giving $P(\text{class}|h_{[CLS]})$ the probabilities of the output vector $h_{[CLS]}$, i.e the whole input sequence, pertaining to a certain *class*.

The training strategies used in this paper are similar to the ones proposed in [48,43]. We use a learning rate $\eta_B = 3e^{-5}$ for the embedding, pooling and encoder layers while a custom learning rate $\eta_C = 1e^{-6}$ for the layers on top of the BERT model. A decay factor $\xi = 1e^{-8}$ is used to reduce gradually the learning rate along the layers, $\eta^l = \xi \cdot \eta^{l-1}$. ADAM optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ and L_2 -weight decay factor of 0.01 is used. The dropout probability

is set at 0.2. Just 5 epochs are enough to find the peak of accuracy with a batch size of 6, the maximum we could use due to memory constraints.

4.3 Image and Text ensemble

In order to get the final enhanced prediction of the combination of both text and image model we use a simple ensemble as in [5].

$$P(\text{class}|\text{out}_{\text{image}}, \text{out}_{\text{text}}) = w_1 \cdot P(\text{class}|h_{[CLS]}) + w_2 \cdot P(\text{class}|FC)$$

$$\text{Predicted Class} = \arg \max_{\text{class}} (P(\text{class}|\text{out}_{\text{image}}, \text{out}_{\text{text}}))$$

In this work $w_1, w_2 = 0.5$ are found optimal. These parameters could be found by a grid search where $\sum_{i=1}^N w_i = 1$, being N the number of models. This procedure shows to be an effective solution when both models have a similar accuracy and it allows us to avoid another training phase [6]. In Figure 2 this whole process is depicted.

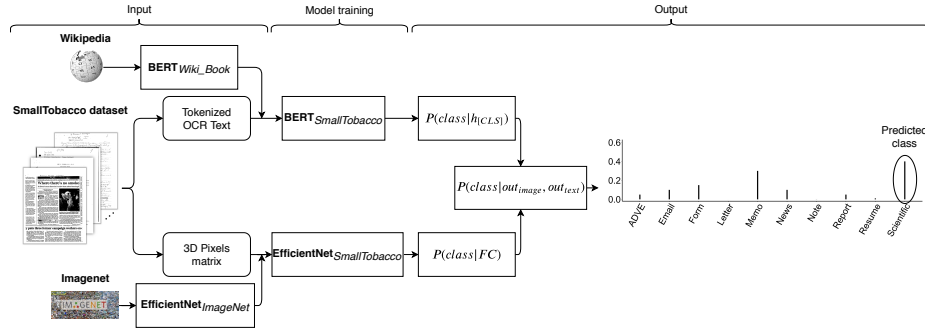


Fig. 2: Pipeline of the proposed multimodal approach.

5 Results

In this section we compare the performance of the different EfficientNets in SmallTobacco and BigTobacco as showed in Table 2 and demonstrate the benefits of the multiple GPU training. Experiments have been carried out using GPUs clusters Power-CTE⁷ of the Barcelona Supercomputing Center - Centro Nacional de Supercomputaci3n⁸, each one composed by: 2 IBM Power9 8335-GTGH at 2.40GHz (20 cores and 4 threads/core), 512GB of main memory distributed in 16 dimms \times 32GB at 2666MHz and 4 GPU NVIDIA V100 (Volta) with 16GB HBM2.

⁷ https://www.bsc.es/support/POWER_CTE-ug.pdf

⁸ <https://www.bsc.es>

The operating system is RedHat Linux 7.4. The models and their training are implemented with PyTorch⁹ version 1.0 running on CUDA 10.1 and using cuDNN 7.6.4.

The only modification done to the images is a resize to 384×384 as explained in Section 4.1 and, in order to avoid overfitting, a shear transformation of an angle $\theta \in [-5^\circ, 5^\circ]$ [46] which is randomly applied in the training phase. No other modifications are used in our experiments. Source code is at <https://javiferran.github.io/document-classification>.

5.1 Evaluation

In order to compare with previous results in SmallTobacco dataset, we divide the dataset following the procedure in [24]. Documents are split in training, test and validation sets, containing 800, 2482 and 200 samples each one. 10 different splits of the dataset are created by randomly sampling from the 3482 documents, so that 100 samples per class are guaranteed between train and validation sets. In the Figure 4 we give the accuracy on SmallTobacco as the median over the 10 dataset splits to compare with previous results. Accuracy on BigTobacco is shown as the one achieved on the test set. BigTobacco dataset used in Section 5.3 is slightly modified, where overlapping documents with SmallTobacco are extracted. Top performing model’s accuracies are written down in Table 2.

5.2 Results on BigTobacco

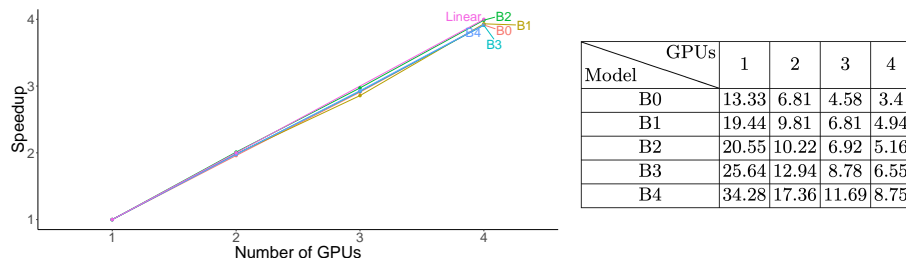


Fig. 3: Left: speedup of the training process when parallelizing. Right: total time (hours) to train each model on different number of GPUs.

We show in Figure 3 the time it takes to train the different networks while using 1, 2, 3 or 4 GPUs in a single node. In order to take advantage of the multiple GPUs we use data parallelism, which consists of placing a copy of the model in each of them. Since every GPU share parameters, it is equivalent to having a single GPU with a larger batch size.

⁹ <https://pytorch.org/>

The time reduction to complete the entire training process with B0 variant is $\approx 61.14\%$ lower when compared with B4 (4 GPUs). Time reduction by using multiple GPUs is clearly showed in the left plot of Figure 3. For instance, EfficientNet-B0 benefits from a $\approx 75.4\%$ time reduction after parallelizing over 4 GPUs. The total training time of the EfficientNets on the different number of GPUs is showed in the right side of Figure 3. The best performing model in BigTobacco dataset is EfficientNet-B4 with 92.31% accuracy in the test set.

5.3 Results on SmallTobacco

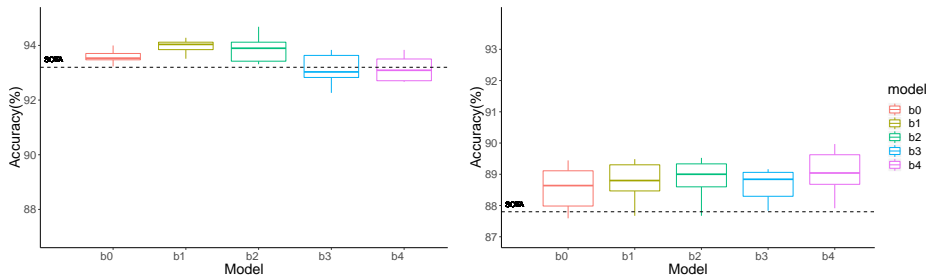


Fig. 4: Accuracy obtained in SmallTobacco by models pre-trained on BigTobacco (Left) and without BigTobacco pre-training (Right). Previous state-of-the-art (SOTA) results are shown with a horizontal dashed line.

Accuracies of the EfficientNets pre-trained on BigTobacco and finally fine-tuned on SmallTobacco are depicted in the left plot of Figure 4. Simpler models perform with less variability between the 10 random splits than the heavier ones. The best performing model is the EfficientNet-B1, achieving a new state-of-art accuracy of 94.04% median over 10 splits.

In this work, we also wanted to test the potential of light EfficientNet models on a small dataset such as SmallTobacco without the use of transfer learning from in-domain dataset, and compared it with the previous state-of-the-art. Results given by our proposed method described in Section 4.3 are shown in the right plot of Figure 4. Although we perform the tests over 10 different random splits to give a wider view of how these models work, in order to compare with *Audebert et al.* [6] we calculate the average over 3 random splits, which gives us a 89.47% accuracy.

Every ensemble model achieves better accuracy than previous results, and again, there is almost no difference between different EfficientNets results.

5.4 Parallel platforms

Single GPU training requires a huge amount of time, especially when dealing with heavy architectures like in the case of the EfficientNet-B4, which takes

almost two days to complete the whole training phase. For this reason, experimenting with several workers is crucial to minimize the amount of time spent on this tasks. We test the same model and training procedure with two of the main used frameworks to train Deep Learning models, PyTorch and Tensorflow¹⁰. In both cases we use their own APIs for making a synchronous distributed training in several GPUs by means of data parallelism, where training on each GPU is done in its own process. We use PyTorch’s DistributedDataParallel and Tensorflow’s `tf.distribute.Strategy` (`tf.distribute.MirroredStrategy`). In both libraries data is loaded from the disk to page-locked memory in each host, and from there to each GPU in a parallel fashion by means of multiple workers. Each GPU is ensured to get a minibatch with non overlapping data. Every GPU has an identical copy of the model and each one does its own forward pass. Finally, NCCL is utilized as a backend to run the all-reduce algorithm to compute the gradients in parallel between GPUs, before updating the model parameters. Since we have not been able to apply the shear transformation efficiently in Tensorflow, we show the results of both frameworks without that preprocess. For this experiment we use the B0, B2 and B4 EfficientNets models. The time it takes to train each model is showed on the left side of Figure 5. PyTorch training is faster and the speedup more linear than in the case of TensorFlow. Some of this difference could be due to the data loading process, which we have not fully optimized in TensorFlow framework.

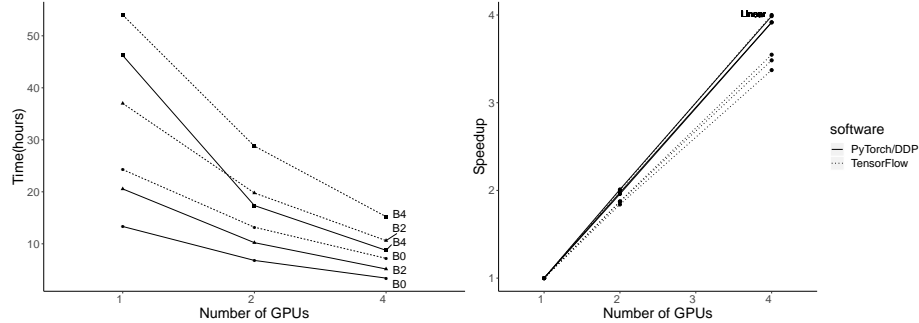


Fig. 5: Left: the time to complete a whole training process. Right: speedup curves of TensorFlow and PyTorch.

6 Conclusion

In this paper we have presented the use of EfficientNets for the Document Image Classification task and their scaling capabilities through several GPUs. By means

¹⁰ <https://www.tensorflow.org/>

of two versions of the Legacy Tobacco Industry Documents, a huge and a small dataset, we demonstrated the training process to obtain high accuracy in both of them. We have compared the different versions of the EfficientNets and raised the state-of-the-art classification accuracy to 92.31% in BigTobacco and 94.04% when fine-tuned in SmallTobacco. We can consider the B0 the best choice when considering limited computational resources. We have also presented an ensemble method by adding the content extracted by OCR. A reduced version of the BERT model is trained and both models predictions are combined to achieve a new state-of-the-art accuracy of 89.47%.

Finally, we have tested the same image models and training procedures in Tensorflow and PyTorch, where we have observed similar speedup values exploiting their libraries for distributed training. We have also tried distributed training in several GPU nodes by means Horovod framework [39], however the stack of software in our IBM Power 9 cluster is still in its early stages and we have not been able to obtain desired results. Nevertheless, future work may focus in testing this approach.

Future work may also evaluate the use of different OCR engines, as we suspect this could have a great impact on the quality of the text model predictions.

With this work we also want to provide to researchers a benchmark in the Document Image Classification task, which can serve as a reference point to effortlessly test parallel systems in both PyTorch and TensorFlow.

7 Acknowledgements

This work was partially supported by the Spanish Ministry of Science and Innovation and the European Regional Development Fund under contract TIN2015-65316-P, by the BSC-CNS Severo Ochoa program SEV-2015-0493, and grant 2017-SGR-1414 by Generalitat de Catalunya and by the research agreement CaixaBank-BSC 2016-2021.

References

1. Afzal, M.Z., Capobianco, S., Malik, M.I., Marinai, S., Breuel, T.M., Dengel, A., Liwicki, M.: Deepdocclassifier: Document classification with deep convolutional neural network. In: ICDAR. p. 1273–1278 (2015)
2. Afzal, M.Z., Kölsch, A., Liwicki, S.A.M.: Cutting the error by half: Investigation of very deep cnn and advanced training strategies for document image classification. In: ICDAR (2017)
3. Aggarwal, C.C.: Neural Networks and Deep Learning: A Textbook. Springer (2018)
4. Akiba, T., Suzuki, S., Fukuda, K.: Extremely large minibatch sgd: Training resnet-50 on imagenet in 15 minutes. arXiv preprint arXiv:1711.04325 (2017)
5. Asim, M.N., Khan, M.U.G., Malik, M.I., Razzaque, K., Dengel, A., Ahmed, S.: Two stream deep network for document image classification. In: ICDAR (2019)
6. Audebert, N., Herold, C., Slimani, K., Vidal, C.: Multimodal deep networks for text and image-based document classification. In: APIA (2019)

7. Baldi, S., Marinai, S., , Soda, G.: Using tree-grammars for training set expansion in page classification. In: ICDAR (2003)
8. Ben-Nun, Hoefler, T.: Demystifying parallel and distributed deep learning: An in-depth concurrency analysis. In: ACM Computing Surveys. vol. 12 (2019)
9. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. In: Trans. Assoc. Comput. Linguist (TACL) (2017)
10. Campos, V., Sastre, F., Yagues, M., Torres, J., i Nieto, X.G.: Scaling a convolutional neural network for classification of adjective noun pairs with tensorflow on gpu clusters. In: CCGRID. pp. 677–682 (2017)
11. Campos, V., Sastre, F., Yagues, M., Torres, M.B.J., i Nieto, X.G.: Distributed training strategies for a computer vision deep learning training algorithm on a distributed gpu cluster. In: ICCS. pp. 315–324 (2017)
12. Chen, S., He, Y., Sun, J., Naoi, S.: Structured document classification by matching local salient features. In: ICPR. pp. 1558–1561 (2012)
13. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: CVPR (2017)
14. Collins-thompson, K., Nickolov, R.: A clustering-based algorithm for automatic document separation. In: SIGIR. p. 1–8 (2002)
15. Csurka, G., Larlus, D., Gordo, A., , Almazan, J.: What is the right way to represent document images? arXiv preprint arXiv:1603.01076 (2016)
16. Das, A., Roy, S., Bhattacharya, U., Parui, S.K.: Document image classification with intra-domain transfer learning and stacked generalization of deep convolutional neural networks. In: ICDAR (2018)
17. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Li, F.F.: Imagenet: a large-scale hierarchical image database. In: CVPR. pp. 248–255 (06 2009)
18. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL (2019)
19. Diligenti, M., Frasconi, P., , Gori, M.: Hidden tree markov models for document image classification. In: Transactions on Pattern Analysis and Machine Intelligence (TPAMI) (2003)
20. Goyal, P., Dollar, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., He, K.: Accurate, large minibatch sgd: Training imagenet in 1 hour. CoRR, vol. abs/1706.02677 (2017)
21. Harley, A.W., Ufkes, A., Derpanis, K.G.: Evaluation of deep convolutional nets for document image classification and retrieval. In: Proc. ICDAR 2015. IEEE. p. 991–995 (2015)
22. He, T., Zhang, Z., Zhang, H., Zhang, Z., Xie, J., Mu Accurate, L.M.S.: Bag of tricks for image classification with convolutional neural networks. arXiv preprint arXiv:1812.01187 (2018)
23. Howard, J., Ruder, S.: Universal language model fine-tuning for text classification. In: Association for Computational Linguistics. vol. 1, p. 328–339 (2018)
24. Jayant, K., Peng, Y., David, D.: Structural similarity for document image classification and retrieval. In: Pattern Recognition Letters. vol. 43, pp. 119–126 (2014)
25. Joulin, A., Grave, E., Bojanowski, P., Mikolov, T.: Bag of tricks for efficient text classification. arXiv preprint arXiv:1607.01759 (2016)
26. Kang, L., Kumar, J., Ye, P., Li, Y., Doermann, D.: Convolutional neural networks for document image classification. In: ICPR. p. 3168–3172 (2014)
27. Kim, Y.: Convolutional neural networks for sentence classification. In: EMNLP (2014)

28. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems* (2012)
29. Kumar, J., Doermann, D.S.: Unsupervised classification of structurally similar document images. In: *ICDAR*. pp. 1225–1229 (2013)
30. Kumar, J., Ye, P., Doermann, D.S.: Learning document structure for retrieval and classification. In: *ICPR*. pp. 653–656 (2012)
31. Lewis, D., Agam, G., Argamon, S., Frieder, O., Grossman, D., Heard., J.: Building a test collection for complex document information processing. In: *SIGIR*. pp. 665–666 (2006)
32. McCann, B., Bradbury, J., Xiong, C., Socher, R.: Learned in translation: Contextualized word vectors. In: *NeurIPS*. pp. 6297–6308 (2017)
33. McCloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: The sequential learning problem. In: *Psychology of learning and motivation*. vol. 24, pp. 109–165 (1989)
34. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. In: *ICLR Workshop Papers* (2013)
35. Noce, L., Gallo, I., Zamberletti, A., Calefati, A.: Embedded textual content for document image classification with convolutional neural networks. In: *Proceedings of the 2016 ACM Symposium on Document Engineering (DocEng '16)* (2016)
36. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: *Proc. of NAACL* (2018)
37. Roy, S., Das, A., Bhattacharya, U.: Generalized stacking of layerwise-trained deep convolutional neural networks for document image classification. In: *23rd International Conference on Pattern Recognition (ICPR)*. p. 1273–1278 (2016)
38. Sandler, M., Howard, A., Menglong, Zhu, Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: *CVPR*. pp. 4510–4520 (2018)
39. Sergeev, A., Balso, M.D.: Horovod: fast and easy distributed deep learning in TensorFlow. *arXiv preprint arXiv:1802.05799* (2018)
40. Shin, C., Doermann, D.S.: Document image retrieval based on layout structural similarity. In: *IPCV*. pp. 606–612 (2006)
41. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556 (2014)
42. Smith, R.: An overview of the tesseract ocr engine. In: *International Conference on Document Analysis and Recognition (ICDAR)* (2007)
43. Sun, C., Qiu, X., Xu, Y., Huang, X.: How to fine-tune bert for text classification? *arXiv preprint arXiv:1905.05583* (2019)
44. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: Mnasnet: Platform-aware neural architecture search for mobile. In: *CVPR* (2019)
45. Tan, M., Le, Q.V.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International Conference on Machine Learning* (2019)
46. Tensmeyer, C., Martinez, T.: Analysis of convolutional neural networks for document image classification. In: *ICDAR* (2017)
47. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: *Advances in Neural Information Processing Systems* 30. p. 6000–6010 (2017)
48. Wang, R., Su, H., Wang, C., Ji, K., Ding, J.: To tune or not to tune? how about the best of both worlds? *arXiv preprint arXiv:1907.05338* (2019)
49. Zhang, Y., Wallace, B.C.: A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. *arXiv preprint arXiv:1510.03820* (2015)